

Assignment 8: Introduction to Random Forest Classifier with Python

Shiva Agarwal

May 16, 2024

Assignment Questions

Question 1: Reading and Exploring the Wine Dataset

1. Read the wine dataset stored in from the location <https://gist.github.com/tijptjik/9408623> and store it as a DataFrame.
2. Assign appropriate column names to the dataset using the given attribute names.
3. Print the first five rows of the DataFrame to examine the dataset.
4. Display the statistical insights of the dataset using the `describe()` function.
5. Visualize the distribution of the 'Alcohol' attribute using a histogram.
6. Visualize the distribution of the 'Ash' attribute using a histogram.
7. Visualize the distribution of the 'Class' attribute using a histogram.
8. Create pair plots to visualize the correlations between selected attributes.
9. Generate a correlation matrix and visualize it using a heatmap.

Question 2: Preparing the Data for Model Training

1. Separate the independent variables (features) and the dependent variable (target) into two separate variables, X and y respectively.
2. Split the dataset into training and testing sets with a test size of 33% and a random state of 42.

Question 3: Building and Evaluating the Random Forest Classifier

1. Create a RandomForestClassifier model.
2. Fit the model to the training data.
3. Predict the target variable using the fitted model.
4. Print the accuracy score of the model.
5. Print the confusion matrix to evaluate the performance of the model.

Solution

Question 1: Reading and Exploring the Wine Dataset

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import confusion_matrix
7 from sklearn.ensemble import RandomForestClassifier
8 from sklearn.metrics import accuracy_score
9
10 # Read the dataset
11 df = pd.read_csv('/Datasets/RandomForests/wine.data')
12
13 # Assign attribute names
14 attributes = ["Class", "Alcohol", "Malic-acid", "Ash", "
15               "Alcalinity-of-ash", "Magnesium",
16               "Total-phenols", "Flavanoids", "
17               "Nonflavanoid-phenols", "Proanthocyanins"
18               "Color-intensity", "Hue", "OD280/OD315-of-
19               diluted-wines", "Proline"]
20 df.columns = attributes
21
22 # Print the first few rows of the DataFrame
23 print("First-few-rows-of-the-DataFrame:")
24 print(df.head())
25
26 # Display statistical insights
27 print("\nStatistical-insights-of-the-dataset:")
28 print(df.describe())
29
30 # Visualize 'Alcohol' using a histogram
31 df.hist('Alcohol')
32 plt.title("Histogram-for-Alcohol")
33 plt.show()
34
35 # Visualize 'Ash' using a histogram
36 df.hist('Ash')
37 plt.title("Histogram-for-Ash")
38 plt.show()
39
40 # Visualize 'Class' using a histogram
41 df.hist('Class')
```

```

39 plt.title("Histogram for Class")
40 plt.show()
41
42 # Create pair plots
43 df_n = df[["Class", "Alcohol", "Malic acid", "Flavanoids",
44           "Nonflavanoid phenols",
45           "Proanthocyanins", "Color intensity", "Hue", "
46           OD280/OD315 of diluted wines",
47           "Proline"]]
48
49 sns.pairplot(df_n, height=4, kind="reg", markers=".")
50 plt.show()
51
52 # Generate a correlation matrix and visualize it using a
53 # heatmap
54 corr = df.corr()
55 cmap = sns.diverging_palette(220, 10, as_cmap=True)
56 sns.heatmap(corr, cmap=cmap, vmax=.3, square=True,
57             linewidths=6, cbar_kws={"shrink": .5})
58 plt.show()
59
60 # Detailed correlation matrix
61 plt.figure(figsize=(12,12))
62 plt.title('Pearson Correlation of Features', y=1.05, size
63           =15)
64 sns.heatmap(df.corr(), linewidths=0.1, vmax=1.0, square=
65             True, cmap=plt.cm.viridis,
66             linecolor='white', annot=True)
67 plt.show()

```

Question 2: Preparing the Data for Model Training

```

1 # Separate features and target variable
2 y = df['Class']
3 X = df.drop(columns=['Class'])
4
5 # Split dataset into training and testing sets
6 X_train, X_test, y_train, y_test = train_test_split(X, y,
7             test_size=0.33, random_state=42)

```

Question 3: Building and Evaluating the Random Forest Classifier

```

1 # Create RandomForestClassifier model

```

```

2 | classifier = RandomForestClassifier(n_jobs=2,
   |     random_state=42)
3 |
4 | # Fit the model to training data
5 | classifier.fit(X_train, y_train)
6 |
7 | # Predict the target variable
8 | y_pred = classifier.predict(X_test)
9 |
10 | # Print the accuracy score
11 | print("Accuracy:", accuracy_score(y_test, y_pred))
12 |
13 | # Print confusion matrix
14 | print("\nConfusion-Matrix:")
15 | print(confusion_matrix(y_test, y_pred))

```